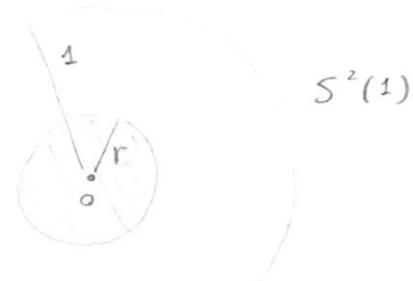


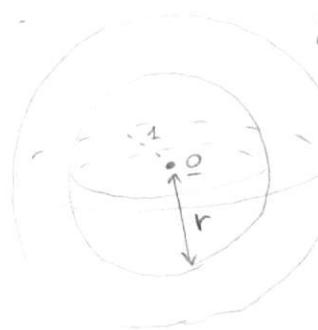
$$p=2: \quad x \sim U[S^2(1)].$$

$$0.1 = \frac{\text{Area}(S^2(r))}{\text{Area}(S^2(1))} = \frac{2\pi r^2}{2\pi} = r^2.$$



$$\Rightarrow r = \sqrt{\frac{1}{10}} = 0.31.$$

$$p=3:$$



$$S^3(1)$$

$$0.1 = \frac{\frac{4}{3}\pi r^3}{\frac{4}{3}\pi} = r^3.$$

$$\Rightarrow r = \left(\frac{1}{10}\right)^{1/3} = 0.46.$$

If you want to capture 10% of your friends, you have to travel 50% of the radius of the universe.

Going on, $\begin{cases} p=100 \Rightarrow r_{0.1} = 0.97 \\ p=1000 \Rightarrow r_{0.1} = 0.99. \end{cases}$

The equation is:

$$r^p = 0.1.$$

When p is big, we will have many isolated coses, very different one from each other.

How do we avoid the curse of dimensionality?

1) Reduce the dimensionality p . We could do it with Principal Component Analysis.

2) Use a parametric model:

$$\text{For example, } f(\mathbf{x}) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p.$$

Everything is reduced to estimating a given small number p of parameters $(\beta_0, \dots, \beta_p)$.

The first approach is data driven.

$$Y = f(X) + \varepsilon.$$

Options:

$$f(x) = \beta_0 + \beta_1 x.$$

$$f(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3.$$

$$f(x) = \beta_0 + \beta_1 x + \dots + \beta_K x^K.$$

We have many different models.

The third one reduces completely the error, i.e. fits perfectly the data. But we see that we have a Variance - Bias tradeoff.

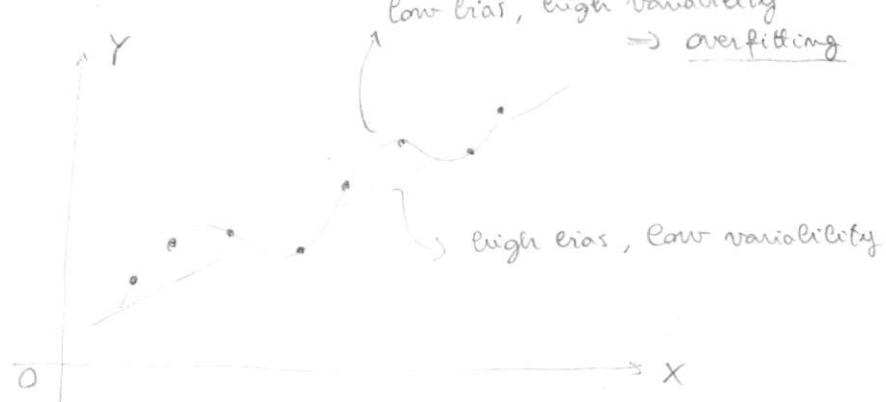
$$\mathbb{E}_{\mathbf{X}} \left[(Y_0 - \hat{f}(x_0))^2 \right] = (f(x_0) - \hat{f}(x_0))^2 + \text{Var}(\varepsilon_0).$$

I want to take the expected value all over the set of datas $\{\mathbf{X}\}$:

$$\begin{aligned} \mathbb{E} \left[\mathbb{E}_{\mathbf{X}} \left[(Y_0 - \hat{f}(x_0))^2 \right] \right] &= \mathbb{E} \left[(f(x_0) - \hat{f}(x_0))^2 \right] + \text{Var}(\varepsilon_0) = \\ &= \mathbb{E} \left[(f(x_0) - \mathbb{E}[\hat{f}(x_0)])^2 \right] + \mathbb{E} \left[(\mathbb{E}[\hat{f}(x_0)] - \hat{f}(x_0))^2 \right] + \text{Var}(\varepsilon_0) = \\ &= \mathbb{E} \left[(f(x_0) - \mathbb{E}[\hat{f}(x_0)])^2 \right] + \mathbb{E} \left[(\mathbb{E}[\hat{f}(x_0)] - \hat{f}(x_0))^2 \right] + 2\mathbb{E}[...] + \text{Var}(\varepsilon_0) = \\ &= \underbrace{(f(x_0) - \mathbb{E}[\hat{f}(x_0)])^2}_{\text{squared bias}} + \underbrace{\text{Var}(\hat{f}(x_0))}_{\text{Variance of the model}} + \text{Var}(\varepsilon_0). \end{aligned}$$

Bias: difference from what you predict and what you observe. (on average)

$$\Rightarrow \boxed{\mathbb{E}[(Y_0 - \hat{f}(x_0))^2] = (f(x_0) - \mathbb{E}[\hat{f}(x_0)])^2 + \text{Var}(\hat{f}(x_0)) + \text{Var}(\varepsilon_0)}.$$



13/3/20

Presentation of Datasets:Data requirements:

- multivariate data;
- two or more groups;

Websites:

Kaggle
Comune di Milano
Regione Lombardia
ISTAT
EUROSTAT
NASA

Project steps:

- 1) Identification of stakeholders (who is interested in your analysis)
- 2) Identification of research questions
- 3) Building the data set from raw data
- 4) Data analysis
- 5) Answer to the research questions

Have in mind somebody asking you some questions (the stakeholders).

At the end, we will present a poster.

Some Possible Datasets:

(Riccardo Peli)

- 1) West Nile Virus Prediction.
- 2) Predicting Molecular Properties.

(Agostino Tati)

- 3) Corona Virus dataset.
- 4) Atlas of Economic Complexity (Harvard University)

Question: predict the spread of Corona Virus in next months.

(Alex Dikkenkampf)

- 5) Estimate the unit sales of Walmart retail goods.
- 6) Factors that help measure how young children learn.
(Marta Fontana)
- 7) Radar Interferometry
(Mara Bernardi)
- 8) Spectroscopic Data for Planetary Investigations.
(Alessandra Menafoglio)
- 9) Earthquake intensity measures.
(Marta Speziori)
- 10) Osteosarcoma
(Daniela Sciacqua)
- 11) Chest Radiography
- 12) Data analysis for Mobility (Agostino Torti) (Marta Galvani)
 - ↳ Ecotra
 - ↳ Train delays
 - ↳ changing network infrastructures
- 13) Safari Njema \Rightarrow "Buen Maggio".
 - Transport Poverty: lack of mobility services
 - ↳ Problem: Informal Mobility.
 - ↳ Solution: 70% of population own a mobile phone.
 - Maputo, Mozambique.
- 14) Identity Recognition via Active Phases.
 - ↳ "Individual Mobility" problem

- n : statistical units
- p : features (or variables)

$$\mathbb{X} = \begin{bmatrix} x_1 & \dots & x_p \\ 1 & \left[\begin{array}{cccc} x_{11} & x_{12} & \dots & x_{1p} \\ \vdots & & & \vdots \\ n & x_{n1} & \dots & x_{np} \end{array} \right] \end{bmatrix}$$

data matrix (data frame)

→ i -th multivariate
observation

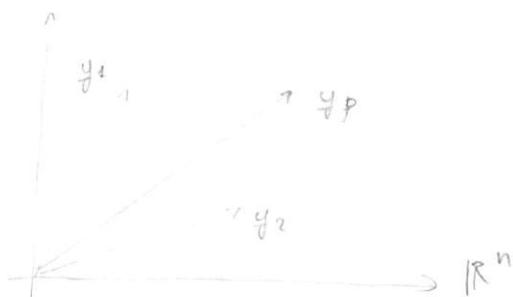
Two uses: by rows and by columns.

- $\underline{x}_i^t = (x_{i1}, \dots, x_{ip}) \in \mathbb{R}^p$, $i = 1, \dots, n$.

$$\underline{y}_1 = (x_{11}, x_{21}, \dots, x_{n1})^t \in \mathbb{R}^n$$

↳ sample from X_1

- $\underline{y}_j = (x_{1j}, \dots, x_{nj})^t \in \mathbb{R}^n$, $j = 1, \dots, p$

↳ sample from X_j • We can take the mean of these samples:

$$\bar{\underline{y}}_j = (x_{1j}, \dots, x_{nj})^t \in \mathbb{R}^n$$

$$\bar{x}_{ij} := \frac{1}{n} \sum_{i=1}^n x_{ij}, \quad j = 1, \dots, p$$

$$\rightarrow \bar{\underline{x}} := \begin{pmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{pmatrix} \quad \text{sample mean vector}$$

• Variance:

$$s_{jj} := \frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2 \quad \text{sample variance of } \underline{x}_j$$

The standard deviation of x_j is $\sqrt{s_{jj}}$.• Covariance between x_k, x_j :

$$s_{kj} := \text{Cov}(\underline{x}_k, \underline{x}_j) = \frac{1}{n} \sum_{i=1}^n (x_{ik} - \bar{x}_k)(x_{ij} - \bar{x}_j), \quad k, j = 1, \dots, p.$$

So that $\text{Cov}(\underline{x}_j, \underline{x}_j) = s_{jj}$.

So we can introduce the covariance matrix (sample):

$$S = \begin{pmatrix} S_{11} & S_{12} & \dots & S_{1p} \\ \vdots & \ddots & \ddots & \vdots \\ & & \ddots & S_{pp} \end{pmatrix}$$

Sample variance / covariance matrix

Properties of S :

i) symmetric.

• Correlation:

$$r_{kj} := \frac{S_{kj}}{\sqrt{S_{kk} S_{jj}}} = \text{Cor}(x_k, x_j), \quad k, j = 1, \dots, p.$$

We can collect all the informations in a matrix:

$$r := \begin{pmatrix} r_{11} & \dots & r_{1p} \\ \vdots & \ddots & \vdots \\ & & r_{pp} \end{pmatrix}$$

Sample correlation matrix

We know that $r_{kj} \in [-1, 1]$.

What do we do with mean and standard deviation?

Example: X_1 : height of an Italian male.

$$\bar{x}_1 = 1.80 \text{ m}, \quad \sqrt{S_{11}} = 0.03 \text{ m}.$$

We know that at least 90% of the population is inside the interval:

$$[\bar{x}_1 - 3\sqrt{S_{11}}, \bar{x}_1 + 3\sqrt{S_{11}}] = [1.74, 1.86].$$

This fact comes from Chebyshev inequality:

$$\Pr [\bar{x}_1 - K\sqrt{S_{11}} \leq X_1 \leq \bar{x}_1 + K\sqrt{S_{11}}] \geq 1 - \frac{1}{K^2}.$$

$$K=3: \quad \Pr [\bar{x}_1 - 3\sqrt{S_{11}} \leq X_1 \leq \bar{x}_1 + 3\sqrt{S_{11}}] \geq \frac{8}{9}.$$

$$K=2: \quad \Pr [\dots] \geq \frac{3}{4}.$$

This without any hypothesis on the distribution of X_1 .

A bit of geometry:

$\underline{v}, \underline{w} \in \mathbb{R}^n$. In \mathbb{R}^n we have an inner product:

$$\langle \underline{v}, \underline{w} \rangle := \underline{v}' \underline{w}.$$

We also have a notion of length:

$$\|\underline{v}\| := \sqrt{\langle \underline{v}, \underline{v} \rangle} = \sqrt{\sum v_i^2}.$$

Moreover, we have the notion of angle:

$$\cos \theta := \frac{\langle \underline{v}, \underline{w} \rangle}{\|\underline{v}\| \|\underline{w}\|} = \frac{\sum v_i w_i}{\sqrt{(\sum v_i^2)(\sum w_i^2)}}$$

Projection of \underline{v} on \underline{w} :

$$\pi_{\underline{v}|\underline{w}} = \pi_{\underline{v}|\mathcal{L}(\underline{w})} = \|\underline{v}\| \cos \theta \cdot \frac{\underline{w}}{\|\underline{w}\|} =$$

$$\mathcal{L}(\underline{w}) = \{ \underline{z} : \underline{z} = c\underline{w}, c \in \mathbb{R} \}.$$

$$= \frac{\|\underline{v}\| \underline{v}' \underline{w}}{\|\underline{v}\| \|\underline{w}\|} \cdot \frac{\underline{w}}{\|\underline{w}\|} = \frac{\underline{v}' \underline{w}}{\|\underline{w}\|^2} \cdot \underline{w} =$$

$$= \frac{\underline{w}' \underline{v}}{\underline{w}' \underline{w}} \cdot \underline{w} = \begin{pmatrix} \underline{w} & \underline{w}' \\ \underline{w}' & \underline{w}' \end{pmatrix} \cdot \underline{v}.$$

(this is a matrix, that projects \underline{v} on $\mathcal{L}(\underline{w})$, called orthogonal projection)

Having introduced this geometry, let's go back to our data sample.

$\mathbb{X} = [y_1, \dots, y_p]$, $y_i \in \mathbb{R}^n$ sample from X_i , $i = 1, \dots, p$.

Take, for instance, y_1 as sample of heights.

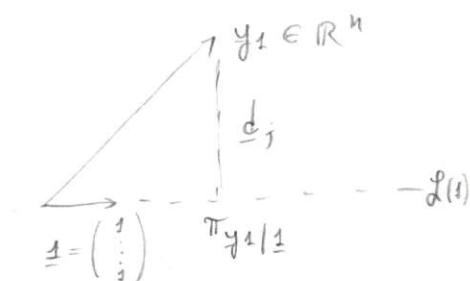
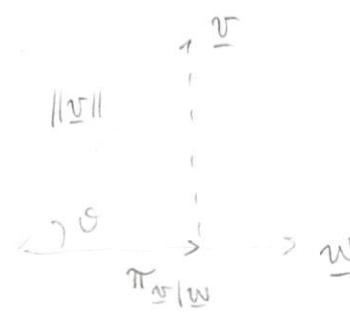
Consider now the $\mathcal{L}(\underline{1})$ space, space of "no statistics" (space where every sample is equal):

$$\underline{v} \in \mathcal{L}(\underline{1}) \Leftrightarrow \underline{v} = c \cdot \underline{1}, c \in \mathbb{R}.$$

Let's project y_1 in $\mathcal{L}(\underline{1})$.

$$\pi_{y_1|\underline{1}} = \frac{\underline{1}' \underline{1}}{\underline{1}' \underline{1}} \cdot y_1 = \frac{\sum_{i=1}^n x_{i1}}{n} \underline{1} = \bar{x}_1 \cdot \underline{1} = \begin{pmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_1 \end{pmatrix}.$$

Hence, geometrically the mean is the projection of the vector in the space with no variability.



anyway, very different vectors can have the same mean, but will have a different deviation \underline{d}_j .

$$\underline{d}_j := \underline{y}_j - \bar{x}_j \cdot \underline{1} = \begin{pmatrix} x_{1j} - \bar{x}_j \\ \vdots \\ x_{nj} - \bar{x}_j \end{pmatrix}. \quad \text{deviation}$$

$$\begin{array}{c} \bar{y}_j \\ \underline{d}_j \\ \bar{x}_j \cdot \underline{1} \end{array} \quad \cancel{\begin{array}{c} \bar{y}_j \\ \underline{d}_j \\ \bar{x}_j \cdot \underline{1} \end{array}} \quad L(\underline{1})$$

The longer is the vector \underline{d}_j , the greater is the norm of \underline{d}_j , the worse is the approximation.

$$\begin{aligned} \|\underline{d}_j\| &= \sqrt{\underline{d}_j \cdot \underline{d}_j} = \sqrt{\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2} \\ &= \sqrt{n} \sqrt{\frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2} = \sqrt{n} \sqrt{s_{jj}}. \end{aligned}$$

Take now two vectors $\underline{y}_j, \underline{y}_k$.

$$\begin{cases} \underline{y}_j = \bar{x}_j \cdot \underline{1} + \underline{d}_j \\ \underline{y}_k = \bar{x}_k \cdot \underline{1} + \underline{d}_k \\ \in L^+(1) \end{cases}$$

$$\begin{array}{c} \underline{d}_j \\ \underline{d}_k \\ \underline{d}_k \end{array} \quad \begin{array}{c} \bar{x}_j \cdot \underline{1} \\ \bar{x}_k \cdot \underline{1} \end{array}$$

Consider the angle between $\underline{d}_j, \underline{d}_k$, say ϑ_{jk} .

- $\vartheta_{jk} = 0$ $\Rightarrow \underline{d}_j \in L(\underline{d}_k) \Rightarrow \underline{d}_j = \beta \underline{d}_k, \beta \in \mathbb{R}$
 $\Rightarrow \underline{y}_j - \bar{x}_j \cdot \underline{1} = \beta (\underline{y}_k - \bar{x}_k \cdot \underline{1})$
 $\Rightarrow \underline{y}_j = \bar{x}_j \cdot \underline{1} + \beta \underline{y}_k - \beta \bar{x}_k \cdot \underline{1}$.

Hence, we have a perfectly linear relationship between the two variables.

- $\vartheta_{jk} = \frac{\pi}{2}$: no information in \underline{d}_k related to \underline{d}_j .
- $\vartheta_{jk} \in (0, \frac{\pi}{2})$: there is some information contained in \underline{d}_k , and some information that is left out.

Therefore, we need to compute the angle ϑ_{jk} in order to know how much one variable is related to the other.

More precisely, we want to measure $\cos \vartheta_{jk}$:

$$\begin{aligned} \cos \vartheta_{jk} &= \frac{\underline{d}_j \cdot \underline{d}_k}{\|\underline{d}_j\| \|\underline{d}_k\|} = \frac{\sum_{i=1}^n (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k)}{\sqrt{\left(\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2\right)\left(\sum_{i=1}^n (x_{ik} - \bar{x}_k)^2\right)}} = \\ &= \frac{s_{jk}}{\sqrt{s_{jj} s_{kk}}} = \text{Cor}(X_j, X_k) = r_{jk} \Rightarrow \text{Therefore it's natural that } r_{jk} \in [-1, 1]. \end{aligned}$$

- $\underline{v}_{jk} = 0 \Rightarrow \cos \underline{v}_{jk} = 1 \Rightarrow r_{jk} = 1 \Rightarrow \underline{d}_j \in \mathcal{L}(\underline{d}_k).$
- $\underline{v}_{jk} = \frac{\pi}{2} \Rightarrow \dots \underline{d}_j \perp \underline{d}_k.$

Let's see now the data matrix from the other perspective:

$$\underline{X} = \begin{bmatrix} \underline{x}_1' \\ \vdots \\ \underline{x}_n' \end{bmatrix}, \quad \underline{x}_i \in \mathbb{R}^p, \quad i=1, \dots, n.$$

\underline{x}_i is the realization of a random vector \underline{X}_i .

(Basic)
General assumption: $\underline{x}_1, \dots, \underline{x}_n$ iid $\sim \underline{X} \in \mathbb{R}^p$.

We need some tools to work with random vectors.

Random Vectors:

Let \underline{X} random vector in \mathbb{R}^p :

$$\underline{X} = \begin{pmatrix} X_1 \\ \vdots \\ X_p \end{pmatrix}, \quad \text{with } X_k, k \in \{1, \dots, p\} \text{ random variable } \in \mathbb{R}.$$

Probability law of \underline{X} : \mathcal{B} : Borel sets of \mathbb{R}^p .

$$v_{\underline{X}} : \mathcal{B} \rightarrow [0, 1].$$

$$v_{\underline{X}}(B) = \Pr[\underline{X} \in B] \quad \forall B \in \mathcal{B}.$$

$$\text{Special cases: } v_{\underline{X}}(B) = \int_B f_{\underline{X}}(t) dt$$

$f_{\underline{X}}$: density of \underline{X} .

$$\text{Mean of } \underline{X}: \quad E[\underline{X}] = \underline{\mu} := \begin{pmatrix} E[X_1] \\ \vdots \\ E[X_p] \end{pmatrix} = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_p \end{pmatrix}$$

Covariance of the components:

$$\sigma_{jk} = E[(X_j - \mu_j)(X_k - \mu_k)] \quad , \quad j=1 \dots p \\ k=1 \dots p$$

$\Sigma := [\sigma_{jk}]$, $p \times p$ matrix.

We can write Σ as:

$$\Sigma = E[(\underline{X} - \underline{\mu})(\underline{X} - \underline{\mu})'].$$

We'll work with combinations (linear) of random vectors.

For this reason, it's comfortable to introduce the matrices:

$$V := \begin{bmatrix} \sigma_{11} & & & \\ & \ddots & & 0 \\ & & \ddots & \\ 0 & & & \sigma_{pp} \end{bmatrix}, \quad \sigma_{ii} > 0$$

$$\tilde{V}^{\frac{1}{2}} = \begin{bmatrix} \sqrt{\sigma_{11}} & & & 0 \\ & \ddots & & \\ 0 & & \sqrt{\sigma_{pp}} & \\ & & & \end{bmatrix} \quad V^{-\frac{1}{2}} = \begin{bmatrix} 1/\sqrt{\sigma_{11}} & & & 0 \\ & \ddots & & \\ 0 & & 1/\sqrt{\sigma_{pp}} & \\ & & & \end{bmatrix} \quad \text{if } \sigma_{ii} > 0.$$

$$g := V^{-\frac{1}{2}} \sum V^{-\frac{1}{2}}$$

correlation matrix ($p \times p$ symmetric)

1) Consider now

$$\underline{\zeta} \in \mathbb{R}^p.$$

$$\Rightarrow \underline{\zeta}' \underline{x} = c_1 x_1 + c_2 x_2 + \dots + c_p x_p \in \mathbb{R}.$$

$$\bullet \quad E[\underline{\zeta}' \underline{x}] = c_1 E[x_1] + \dots + c_p E[x_p] = \underline{\zeta}' \underline{\mu}.$$

$$\bullet \quad \text{Var}(\underline{\zeta}' \underline{x}) = ?$$

Exercise: $\text{Var}(c_1 x_1 + c_2 x_2) = c_1^2 \text{Var}(x_1) + c_2^2 \text{Var}(x_2) + 2c_1 c_2 \text{Cov}(x_1, x_2).$

\Rightarrow I can reduce the uncertainty if $\text{Cov}(x_1, x_2) \leq 0$.

In general,

$$\text{Var}(\underline{\zeta}' \underline{x}) = \underline{\zeta}' \Sigma \underline{\zeta}.$$

2) If we want to take many linear combinations, consider:

C : $k \times p$ matrix of constants.

$$\text{Consider } C\underline{x}. \quad C = \begin{pmatrix} \underline{\zeta}_1' \\ \vdots \\ \underline{\zeta}_k' \end{pmatrix}, \quad \underline{\zeta}_i \in \mathbb{R}^p. \quad \Rightarrow C\underline{x} = \begin{pmatrix} \underline{\zeta}_1' \underline{x} \\ \vdots \\ \underline{\zeta}_k' \underline{x} \end{pmatrix}.$$

$$\bullet \quad E[C\underline{x}] = C\underline{\mu}.$$

$$\bullet \quad \text{Cov}(C\underline{x}) = C X C'$$

Precise case: $C = \underline{\zeta}' = (c_1 \dots c_p)^T \quad 1 \times p.$

$$\text{Cov}(C\underline{x}) = \text{Var}(\underline{\zeta}' \underline{x}) = C \Sigma C' = \underline{\zeta}' \Sigma \underline{\zeta}.$$

Conclusion:

Data: \underline{X} $m \times p$ $\xrightarrow{\text{inference}}$ Model: \underline{X} random vector.

We use the data to estimate $\underline{\mu}, \Sigma$.

• Estimate for μ :

$$\mu = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_p \end{pmatrix} \simeq \begin{pmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{pmatrix} \leftarrow \text{generated by data}$$

But why does it work? We can see $\bar{x} = \begin{pmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{pmatrix}$ as a realization of the object $\frac{1}{n} \sum_{i=1}^n \underline{x}_{i \cdot} =: \bar{X}$. (the estimator of μ).

Notation: $\begin{cases} \text{estimator: algorithm.} \\ \text{estimate: product of the algorithm.} \end{cases}$

Proposition: Let $\underline{x}_1, \dots, \underline{x}_n$ iid $\sim \underline{x}$, $E[\underline{x}] = \mu$, $\text{cov}(\underline{x}) = \Sigma$. Then:

$$1. E[\bar{X}] = \mu \quad (\text{i.e. } \bar{X} \text{ is unbiased}).$$

$$2. \text{cov}(\bar{X}) = \frac{1}{n} \Sigma.$$

Proof:

$$E[\bar{X}] = E \left[\begin{pmatrix} \frac{1}{n} \sum_{i=1}^n x_{i1} \\ \vdots \\ \frac{1}{n} \sum_{i=1}^n x_{ip} \end{pmatrix} \right] = \mu.$$

$$\begin{aligned} \text{cov}(\bar{X}) &= E[(\bar{X} - \mu)(\bar{X} - \mu)'] = E \left[\left(\frac{1}{n} \sum x_i - \mu \right) \left(\frac{1}{n} \sum x_i - \mu \right)' \right] = \\ &= \frac{1}{n} \sum_{i=1}^n E[(x_i - \mu)(x_i - \mu)'] \\ &= E \left[\frac{1}{n} \sum (x_i - \mu) \left(\frac{1}{n} \sum (x_i - \mu) \right)' \right] = \\ &= \frac{1}{n^2} E \left[\sum_{i=1}^n \sum_{k=1}^n (x_i - \mu)(x_k - \mu)' \right] = \\ &= \frac{1}{n^2} \sum_{i, k=1}^n E[(x_i - \mu)(x_k - \mu)']. \end{aligned}$$

$$\text{Now, } E[(x_i - \mu)(x_k - \mu)] = \begin{cases} \Sigma, & i=k, \\ 0, & i \neq k \text{ (independence)} \end{cases}$$

$$\Rightarrow \text{cov}(\bar{X}) = \frac{1}{n^2} \sum_{i=1}^n (\Sigma) = \frac{n}{n^2} \Sigma = \frac{1}{n} \Sigma. \quad (\text{end})$$

• Estimate for Σ :

Exercise: let $S := \frac{1}{n} \sum (x_i - \bar{x})(x_i - \bar{x})'$, estimator of Σ .

Prove that $E[S] = \frac{n-1}{n} \Sigma$, i.e. S is biased.

\Rightarrow an unbiased estimator of Σ is $\frac{1}{n-1} S$.

$$\underline{X} = \begin{bmatrix} \underline{x}_1 \\ \vdots \\ \underline{x}_n \end{bmatrix} = [y_1 \ \dots \ y_p] \quad \begin{array}{l} \underline{x}_i \in \mathbb{R}^p \\ \text{data matrix} \end{array}$$

\underline{x}_i is the realization of a random vector \underline{X}_i in \mathbb{R}^p .

We will assume that $\underline{x}_1, \dots, \underline{x}_n$ are i.i.d. $\sim \underline{X}$, with:

$$E[\underline{X}] = \mu, \text{cov}(\underline{X}) = \Sigma \quad (\text{unknown}).$$

In order to estimate these parameters, we introduce:

$$\bar{\underline{X}} = \frac{1}{n} \sum_{i=1}^n \underline{x}_i \quad \text{random vector in } \mathbb{R}^p, \text{ called sample mean.}$$

The realization of $\bar{\underline{X}}$ is $\bar{\underline{x}} = \begin{pmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{pmatrix}$.

$$\underline{S}' = \frac{1}{n} \sum_{i=1}^n (\underline{x}_i - \bar{\underline{x}})(\underline{x}_i - \bar{\underline{x}})' \quad \text{is a } p \times p \text{ random matrix}$$

$$\text{The realization of } \underline{S}' \text{ is } \underline{S} = \begin{bmatrix} s_{11} & \dots & s_{1p} \\ \vdots & & \vdots \\ s_{p1} & \dots & s_{pp} \end{bmatrix}, \text{ sample cov. matrix.}$$

What we've seen yesterday is that:

Proposition:

1). $E[\bar{\underline{X}}] = \mu$, i.e. $\bar{\underline{X}}$ is unbiased for μ .

$\text{Cov}(\bar{\underline{X}}) = \frac{1}{n} \Sigma$, i.e. the covariance of $\bar{\underline{X}}$ is smaller than Σ .

2). $E[\underline{S}'] = \frac{n-1}{n} \Sigma$, i.e. \underline{S}' is a biased estimator.

A simple corollary is: $E\left[\frac{n}{n-1} \underline{S}'\right] = \Sigma$.

So $\frac{n}{n-1} \underline{S}'$ is an unbiased estimator for Σ .

$$\frac{m}{m-1} \underline{S}' = \frac{m}{m-1} \frac{1}{n} \sum_{i=1}^n (\underline{x}_i - \bar{\underline{x}})(\underline{x}_i - \bar{\underline{x}})' = \frac{1}{n-1} \sum_{i=1}^n (\underline{x}_i - \bar{\underline{x}})'(\underline{x}_i - \bar{\underline{x}})'.$$

From now on, we will denote:

$$\underline{S}' = \frac{1}{n-1} \sum_{i=1}^n (\underline{x}_i - \bar{\underline{x}})(\underline{x}_i - \bar{\underline{x}})'.$$

Where we'll need the "dd \underline{S}' ", i.e. the covariance of the sample, we will denote it with S_n .

Today we will focus on \underline{S}' . Some alternative notations for \underline{S}' will be found. Remind the definition of deviation vector d_j :

$$d_j = y_j - \pi_{y_j \perp} = y_j - \frac{1 \cdot 1'}{1' 1} y_j = (I - \underbrace{\frac{1 \cdot 1'}{1' 1} I}_{\text{orthogonal projector on } \mathbb{L}^\perp(1)}) y_j.$$

orthogonal projector on $\mathbb{L}^\perp(1)$.

The deviation matrix is defined by:

$$d := [d_1 \ d_2 \ \dots \ d_p] \text{, } n \times p \text{ matrix.}$$

We can see that:

$$d = (I - \frac{11'}{1'1}) X.$$

Moreover,

$$\begin{aligned} S' &= \begin{bmatrix} s_{11} & s_{12} & \dots & s_{1p} \\ & \ddots & \ddots & \vdots \\ & & \ddots & s_{pp} \end{bmatrix} = \frac{1}{n-1} \begin{bmatrix} \sum (x_{i1} - \bar{x}_1)^2 & \sum (x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2) \dots \\ \vdots & \vdots \\ \sum (x_{ip} - \bar{x}_p)^2 & \sum (x_{ip} - \bar{x}_p)(x_{i1} - \bar{x}_1) \dots \end{bmatrix} = \\ &= \frac{1}{n-1} \begin{bmatrix} d_1'd_1 & d_1'd_2 & \dots & d_1'd_p \\ & \vdots & & \vdots \\ & d_p'd_1 & d_p'd_2 & \dots & d_p'd_p \end{bmatrix} = \frac{1}{n-1} d'd = \\ &= \frac{1}{n-1} X'(I - \frac{11'}{1'1})'(I - \frac{11'}{1'1})X. \end{aligned}$$

Given that $I - \frac{11'}{1'1}$ is an orthogonal projection,

$$(I - \frac{11'}{1'1})'(I - \frac{11'}{1'1}) = (I - \frac{11'}{1'1}).$$

$$\Rightarrow S' = \frac{1}{n-1} X'(I - \frac{11'}{1'1})X.$$

Variability in a "multivariate" sense:

How do we capture the variability of multi-variate data? We have two options:

1) Generalized variance: $\text{Det}(S)$.

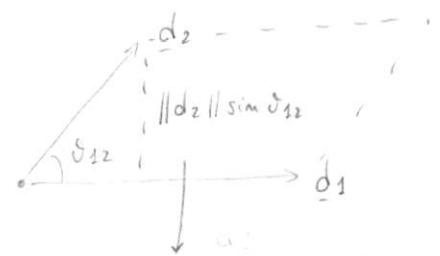
2) Total Variance: $\text{tr}(S)$.

Example: consider $p=2$. Then:

$$\begin{aligned} S' &= \frac{1}{n-1} d'd = \frac{1}{n-1} \begin{bmatrix} d_1'd_1 & d_1'd_2 \\ d_2'd_1 & d_2'd_2 \end{bmatrix} = \\ &= \frac{1}{n-1} \begin{bmatrix} \|d_1\|^2 & \|d_1\| \|d_2\| \cos \vartheta_{12} \\ \|d_1\| \|d_2\| \cos \vartheta_{12} & \|d_2\|^2 \end{bmatrix}. \end{aligned}$$

The generalized variance is:

$$\begin{aligned} \text{Det}(S') &= \frac{\|d_1\|^2 \|d_2\|^2 - \|d_1\|^2 \|d_2\|^2 \cos^2 \vartheta_{12}}{(n-1)^2} = \\ &= \frac{\|d_1\|^2 \|d_2\|^2 (1 - \cos^2 \vartheta_{12})}{(n-1)^2} = \frac{\|d_1\|^2 \|d_2\|^2 \sin^2 \vartheta_{12}}{(n-1)^2} = \\ &\propto (\text{Area of parallelogram } \frac{(n-1)^2}{(\underline{d}_1, \underline{d}_2)})^2. \end{aligned}$$



generalized variance:
proportional to the area²
of the parallelogram

Then, $\det(S') \uparrow \Leftrightarrow$ Area parallelogram ↑.

Moreover, $\det(S') = 0 \Rightarrow$ either $\begin{cases} \|\underline{d}_1\| \text{ or } \|\underline{d}_2\| = 0, \\ \underline{d}_{12} = 0 \end{cases}$

The total variance, instead, is:

$$\text{tr}(S) = (\|\underline{d}_1\|^2 + \|\underline{d}_2\|^2) \frac{1}{n-1}.$$

This holds for p in general:

$$\begin{cases} \det(S) \propto \text{Vol}^2(\text{parallelepiped } (\underline{d}_1, \dots, \underline{d}_p)). \\ \text{trace}(S) \propto \|\underline{d}_1\|^2 + \|\underline{d}_2\|^2 + \dots + \|\underline{d}_p\|^2. \end{cases}$$

We need both these variances in order to describe the variability of the multivariate data.

Proposition:

$\det(S) = 0 \Leftrightarrow \underline{d}_1, \dots, \underline{d}_p$ are linearly dependent,
i.e. $\exists \underline{c} \neq \underline{0}$ s.t. $\underline{d}\underline{c} = \underline{0}$ (i.e. $c_1 \underline{d}_1 + \dots + c_p \underline{d}_p = \underline{0}$).

Proof:

\Leftarrow Suppose $\underline{d}_1, \dots, \underline{d}_p$ are linearly dependent, i.e.
 $\exists \underline{c} \in \mathbb{R}^p$ s.t. $\underline{d}\underline{c} = \underline{0}$.

Then: $S = \frac{1}{n-1} \underline{d}' \underline{d} \Rightarrow S\underline{c} = \frac{1}{n-1} \underline{d}' \underbrace{\underline{d}\underline{c}}_{= \underline{0}} = \underline{0}$

\Rightarrow the columns of S are linearly dependent $\Rightarrow \det(S) = 0$.

\Rightarrow Assume $\det(S) = 0$. Then $\exists \underline{c} \neq \underline{0}$ s.t. $S\underline{c} = \underline{0}$.

$$\Rightarrow \underline{c}' S \underline{c} = \underline{0} \Rightarrow \frac{1}{n-1} \underline{c}' \underline{d}' \underline{d} \underline{c} = \underline{0}, \text{ or } \frac{1}{n-1} \|\underline{d}\underline{c}\|^2 = \underline{0}.$$

By property of norm, $\|\underline{d}\underline{c}\| = 0 \Rightarrow \underline{d}\underline{c} = \underline{0}$, i.e.

$\underline{d}_1, \dots, \underline{d}_p$ are linearly dependent.

(end)

But what does it mean that $\underline{d}_1, \dots, \underline{d}_p$ are linearly dependent?

We have seen that $\exists c_1 \underline{d}_1 + \dots + c_p \underline{d}_p = \underline{0}$, $\underline{c} \neq \underline{0}$.

W.l.o.g., $c_1 \neq 0$. Then:

$$\underline{d}_1 = - \sum_{i=2}^p \frac{c_i}{c_1} \underline{d}_i$$

$$\Rightarrow y_1 = \bar{x}_1 \cdot 1 - \sum_{i=2}^p \frac{c_i}{c_1} (y_i - \bar{x}_i \cdot 1).$$

Hence there is no new information from variable x_1 : we have a perfect linear relationship between the first variable and the others.

Proposition: Let X $m \times p$ data matrix.

If $p \geq n$, then $\det(S) = 0$.

Proof: $d = [d_1, \dots, d_p]$, $d_i \in \mathbb{R}^n$, but $d_i \in L^+(1)$, $i=1 \dots n$

$\dim(L^+(1)) = n-1$: $[d_i]_{i=1}^p$ live in a space of dimension $n-1$.
($n-1$ degrees of freedom)

But $p \geq n \Rightarrow d_1, \dots, d_p$ are linearly dependent.

$\Rightarrow \det(S) = 0$.

$$\begin{array}{cc} X & Y \\ \left[\begin{array}{c} x_1 \\ \vdots \\ x_n \end{array} \right] & \left[\begin{array}{c} y_1 \\ \vdots \\ y_n \end{array} \right] \end{array} \quad (\text{end})$$

Example: suppose our raw data is:

Note that the data matrix X doesn't have to be equal to the raw data. For example:

$$X = \left[\begin{array}{cccccc} x & x^2 & \dots & x^p & y \\ x_1 & x_1^2 & \dots & x_1^p & y_1 \\ x_2 & x_2^2 & \dots & x_2^p & y_2 \\ \vdots & \vdots & & \vdots & \vdots \\ x_n & x_n^2 & \dots & x_n^p & y_n \end{array} \right].$$

Taking $p = n-1$, then the number of columns of X is $p+1 = n$.
 $\Rightarrow \det(S) = 0$. Assume that x_n, x_n^2, \dots, x_n^p are linearly independent.

Then: $y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \dots + \beta_p x_i^p$.

So we have a perfect fit \Rightarrow overfitting! Remember the bias-variance tradeoff.

This example tells us we'll always be able to find a perfect fit for the data, but that's not what we want.

Spectral Decomposition of S :

We have seen that a necessary condition for $\det(S) \neq 0$ is $n \geq p+1$.

We know that S is a $p \times p$ symmetric matrix, and $s_{ij} \in \mathbb{R} \quad \forall i, j = 1, \dots, p$.
From linear algebra, we do know that:

$\exists \lambda_1, \dots, \lambda_p \in \mathbb{R}$,

$e_1, \dots, e_p \in \mathbb{R}^p$ s.t. $e_i^T e_j = \begin{cases} 0, & i \neq j \\ 1, & i = j \end{cases}$

$$S = \sum_{i=1}^p \lambda_i e_i e_i^T$$

spectral decomposition of S

In fact, (λ_i, e_i) is the eigenvalue - eigenvector couple of S .

This decomposition is the core of dimensional reduction.

Let: $P := [\underline{e}_1, \dots, \underline{e}_p]$, $\Lambda = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_p \end{bmatrix}$.

Since $\underline{e}_i^T \underline{e}_i = \lambda_i \underline{e}_i^T \underline{e}_i \quad \forall i=1, \dots, p$, then:

$$\underline{S} = P \Lambda P'$$

Spectral decomposition of \underline{S}
(matrix form)

Moreover,

$$\left\{ \begin{array}{l} \det(\underline{S}) = \prod_{i=1}^p \lambda_i \quad \text{generalized variance} \\ \text{tr}(\underline{S}) = \sum_{i=1}^p \lambda_i \quad \text{total variance} \end{array} \right.$$

Proposition: \underline{S} is positive semi-definite.

If $\det(\underline{S}) \neq 0$, then \underline{S} is positive definite.

Proof: need to prove that $\underline{c}' \underline{S} \underline{c} \geq 0 \quad \forall \underline{c} \in \mathbb{R}^p$.

$$\text{But } \underline{c}' \underline{S} \underline{c} = \frac{1}{n-1} \underline{c}' \underline{d}' \underline{d} \underline{c} = \frac{1}{n-1} \|\underline{d} \underline{c}\|^2 \geq 0.$$

Suppose that $\exists \underline{c} \neq \underline{0}$ s.t. $\underline{c}' \underline{S} \underline{c} = 0 \Rightarrow \|\underline{d} \underline{c}\| = 0$.

$$\Rightarrow \underline{d} \underline{c} = \underline{0} \Rightarrow d_1, \dots, d_p \text{ are lin. dep.} \xrightarrow{\text{(prop. before)}} \det(\underline{S}) = 0.$$

So if $\det(\underline{S}) \neq 0$, then \underline{S} is positive definite ($\det(\underline{S}) > 0$).

Mahalanobis Distance and its geometry:

(end)

Notation: from now on, writing $\underline{S} = \sum_{i=1}^p \lambda_i \underline{e}_i \underline{e}_i'$, we assume that:

$$\boxed{\lambda_1 > \lambda_2 > \dots > \lambda_p \geq 0} \quad \text{ordering eigenvalues of } \underline{S}$$

$$\begin{array}{ccc} | & | & | \\ \underline{e}_1 & \underline{e}_2 & \underline{e}_p \end{array}$$

Assume that $\det(\underline{S}) \neq 0 \Rightarrow \lambda_1 \geq \lambda_p > 0$.

So we have:

$$\underline{S} = \sum_{i=1}^p \lambda_i \underline{e}_i \underline{e}_i'$$

$$\Rightarrow \boxed{\underline{S}^{-1} = \sum_{i=1}^p \frac{1}{\lambda_i} \underline{e}_i \underline{e}_i'}$$

Then, \underline{S}^{-1} induces a metric on \mathbb{R}^p . Indeed, let $d_{S^{-1}}$ be defined as:

$$\forall \underline{x}, \underline{y} \in \mathbb{R}^p: \boxed{d_{S^{-1}}^2(\underline{x}, \underline{y}) := (\underline{x} - \underline{y})' \underline{S}^{-1} (\underline{x} - \underline{y})}$$

Mahalanobis distance

$d_{S^{-1}}$ is a distance in \mathbb{R}^p [prove it!].

Consider now the set:

$$\begin{aligned}\mathcal{E}_{r^2, S^{-1}}(\bar{x}) &:= \left\{ \underline{x} \in \mathbb{R}^p : d_{S^{-1}}^2(\underline{x}, \bar{x}) \leq r^2 \right\} \\ &= \left\{ \underline{x} \in \mathbb{R}^p : (\underline{x} - \bar{x})^T S^{-1}(\underline{x} - \bar{x}) \leq r^2 \right\}.\end{aligned}$$

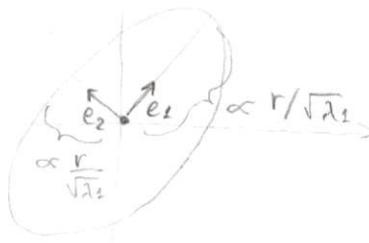
What are the characteristics of this set?

Consider, in general, B positive definite of the form:

$$B = \sum_{i=1}^p \lambda_i e_i e_i^T.$$

Consider the points $\underline{x}^T B \underline{x} = r^2$.

\Rightarrow These points are an ellipse.



In our case, $\mathcal{E}_{r^2, S^{-1}}(\bar{x})$ is the following ellipse:

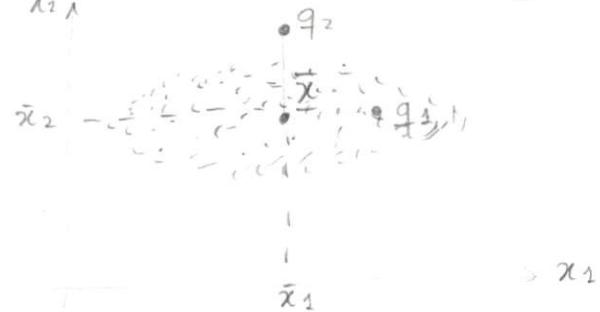
The axes of $\mathcal{E}_{S^{-1}}(\bar{x}, r)$ are proportional to the squared root of the eigenvalues of S .



The volume of $\mathcal{E}_{r^2, S^{-1}}(\bar{x})$ is:

$$\text{Vol}(\mathcal{E}_{r^2, S^{-1}}(\bar{x})) = K_p r^p \sqrt{\prod_{i=1}^p \lambda_i} = K_p r^p \underbrace{\sqrt{\text{Det}(S)}}_{\text{generalized variance}}.$$

Example: Let $p=2$, $S = \begin{bmatrix} s_{11} & 0 \\ 0 & s_{22} \end{bmatrix}$, $s_{11} > s_{22}$.



Take q_1, q_2 s.t: $q_1 = \begin{pmatrix} q_1 \\ \bar{x}_1 \end{pmatrix}$, $q_2 = \begin{pmatrix} \bar{x}_2 \\ q_2 \end{pmatrix}$,

and $d_{\text{eucl}}(\bar{x}, q_1) = d_{\text{eucl}}(\bar{x}, q_2)$.

Anyway, q_1 is a "normal point" with respect to the cloud, while q_2 is an extraordinary point.

We want to deal with standardized data:

$$d_{\text{eucl}}(\text{std}(q_1), \text{std}(\bar{x})) = \frac{|q_1 - \bar{x}_1|}{\sqrt{s_{11}}}.$$

$$d_{\text{eucl}}(\text{std}(q_2), \text{std}(\bar{x})) = \frac{|q_2 - \bar{x}_2|}{\sqrt{s_{22}}}.$$

Given that $|q_1 - \bar{x}_1| = |q_2 - \bar{x}_2|$ and $s_{11} > s_{22}$, then:

$$d_{\text{eucl}}(\text{std}(q_1), \text{std}(\bar{x})) < d_{\text{eucl}}(\text{std}(q_2), \text{std}(\bar{x})).$$

In general, taking $\mathbf{q} = \begin{pmatrix} q_1 \\ q_2 \end{pmatrix}$, we have:

$$\text{d}_{\text{std}}(\text{std}(\mathbf{q}), \text{std}(\bar{\mathbf{x}})) = \sqrt{\frac{(q_1 - \bar{x}_1)^2}{S_{11}} + \frac{(q_2 - \bar{x}_2)^2}{S_{22}}} = \sqrt{(\mathbf{q} - \bar{\mathbf{x}})^T S^{-1} (\mathbf{q} - \bar{\mathbf{x}})} = \text{d}_{S^{-1}}(\mathbf{q}, \bar{\mathbf{x}}).$$

Taking a more general shape of $S' = \begin{bmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{bmatrix}$, we have an inclined ellipse.



Recap:

19/3/20

- \mathbb{X} : data matrix $n \times p$,
- $\bar{\mathbf{x}}$: sample mean $\in \mathbb{R}^p$,
- S' : sample covariance $p \times p$.

If $\det(S) > 0$, $d_{S^{-1}}: \mathbb{R}^p \times \mathbb{R}^p \rightarrow [0, \infty]$.

$$d_{S^{-1}}^2 := (\mathbf{x} - \bar{\mathbf{x}})^T S^{-1} (\mathbf{x} - \bar{\mathbf{x}})$$

Mahalanobis' distance

Example: $S = \begin{bmatrix} S_{11} & 0 \\ 0 & S_{22} \end{bmatrix}$.
($p=2$)

$$\mathcal{E}_r(\bar{\mathbf{x}}) = \{ \mathbf{x} \in \mathbb{R}^p : (\mathbf{x} - \bar{\mathbf{x}})^T S^{-1} (\mathbf{x} - \bar{\mathbf{x}}) \leq r^2 \}.$$



↳ ellipse of "radius" r centered in $\bar{\mathbf{x}}$.

$$\text{Area}(\mathcal{E}_r(\bar{\mathbf{x}})) \propto r^2 \sqrt{\det(S)} = r^2 \sqrt{S_{11} \cdot S_{22}}.$$

The larger the variability $\det(S)$, the larger this neighborhood.

If $\det(S)$ or $\text{tr}(S) \uparrow \Rightarrow \text{Area}(\mathcal{E}_r(\bar{\mathbf{x}})) \uparrow$ (with r fixed)

One could say that it's a special case. Is this really special?

• it's easy to extend to $p > 2$ (use your imagination).

• what if S is not diagonal? Let $p \geq 1$, any S : $\det(S) > 0$.

(In the figure, we have positive correlation between X_1, X_2)



Consider the spectral decomposition of S :

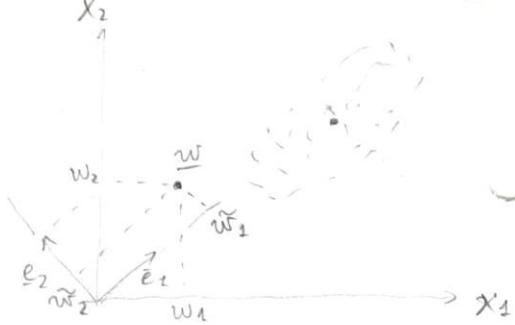
$$S = \sum_{i=1}^p \lambda_i e_i e_i^T, \quad (\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p > 0).$$



Or, in other words:

$$S = P \Lambda P'$$

We'll introduce a new orthonormal system made by the eigenvectors of S .



Old system:

$$\underline{w} = \begin{pmatrix} w_1 \\ \vdots \\ w_p \end{pmatrix}$$

New system:

$$\tilde{\underline{w}} = \begin{pmatrix} e_1' \underline{w} \\ \vdots \\ e_p' \underline{w} \end{pmatrix} = P' \underline{w}$$

Since the length doesn't change:

$$\|\underline{w}\|^2 = \underline{w}' \underline{w} = \tilde{\underline{w}}' \tilde{\underline{w}} = \|\tilde{\underline{w}}\|^2 = \underline{w}' P' P \underline{w}.$$

We can do the same with our data:

$$\mathbb{X} = \begin{bmatrix} \underline{x}_1' \\ \vdots \\ \underline{x}_n' \end{bmatrix} \quad \tilde{\mathbb{X}} = \begin{bmatrix} (P' \underline{x}_1)' \\ \vdots \\ (P' \underline{x}_n)' \end{bmatrix} = \begin{bmatrix} \underline{x}_1' P \\ \vdots \\ \underline{x}_n' P \end{bmatrix} = \mathbb{X} P.$$

Let's analyze the covariance matrix:

$$\begin{aligned} S &= \frac{1}{n-1} \mathbb{X}' \left(I - \frac{\mathbb{1} \mathbb{1}'}{\mathbb{1}' \mathbb{1}} \right) \mathbb{X}, \quad \tilde{S} = \frac{1}{n-1} \tilde{\mathbb{X}}' \left(I - \frac{\mathbb{1} \mathbb{1}'}{\mathbb{1}' \mathbb{1}} \right) \tilde{\mathbb{X}} \\ &\quad \vdots \\ \Rightarrow \tilde{S} &= \Lambda = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_p \end{pmatrix}. \quad \begin{array}{l} \tilde{S} = \frac{1}{n-1} \tilde{\mathbb{X}}' \left(I - \frac{\mathbb{1} \mathbb{1}'}{\mathbb{1}' \mathbb{1}} \right) \tilde{\mathbb{X}} \\ = \frac{1}{n-1} P' \mathbb{X}' \left(I - \frac{\mathbb{1} \mathbb{1}'}{\mathbb{1}' \mathbb{1}} \right) \mathbb{X} P \\ = P' S P = P' P \Lambda P' P \\ = \Lambda. \end{array} \end{aligned}$$

In this new system, the covariances are zero!

Take-home message: there is always a reference system for which the coordinates are uncorrelated.

"We need to look at the data with particular glasses".

Let's now analyze the generalized variance:

$$\text{Det}(\tilde{S}) = \text{Det}(\Lambda) = \prod_{i=1}^p \lambda_i = \text{Det}(S).$$

Some thing goes for total variance:

$$\text{tr}(\tilde{S}) = \text{tr}(\Lambda) = \sum_{i=1}^p \lambda_i = \text{tr}(S).$$

The total variability doesn't depend on the reference system.

Remark:

$$P'P = \begin{pmatrix} e_1 \\ \vdots \\ e_p \end{pmatrix} (e_1 \dots e_p) = \begin{pmatrix} 1 & e_1^T e_1 & e_1^T e_2 & \dots & e_1^T e_p \\ \vdots & \ddots & \ddots & \ddots & 0 \\ e_p^T e_1 & e_p^T e_2 & \dots & e_p^T e_p & 1 \end{pmatrix} = I.$$

Example: (former point). Correlation of the Milky Way.

↳ correlation is in the eyes of the beholder, not in the system itself!

We spot a presence of a 1D linear space: it is the first principal component!

Principal Component Analysis (PCA):

Let \underline{X} random vector in \mathbb{R}^p , $E[\underline{X}] = \mu$, $\text{cov}[\underline{X}] = \Sigma$.

Take $\underline{a} \in \mathbb{R}^p$: $\underline{a}' \underline{X} = a_1 X_1 + \dots + a_p X_p$.

Problem: find \underline{a} s.t. $\text{Var}(\underline{a}' \underline{X})$ is maximum.

Is the problem well-posed?

Suppose that the solution is $\underline{a}' = (1, 0, \dots, 0)'$, s.t. $\max \text{var} = \underline{a}' \underline{X} = 1 = a_1 X_1$.

In this way, we could take:

$$10 \cdot \underline{a} = (100, 0, \dots, 0)$$

$$\Rightarrow \text{Var}(10 \cdot \underline{a}' \underline{X}) = 100 \text{Var}(\underline{a}' \underline{X}).$$

So \underline{a} cannot be the maximum!

Therefore, the problem is not well-posed. Let's refine it a little bit.

Problem: find \underline{a} s.t. $\|\underline{a}\| = 1$ and $\text{Var}(\underline{a}' \underline{X})$ is max.

Formally:

$$\max_{\underline{a} \in \mathbb{R}^p: \|\underline{a}\|=1} \text{Var}(\underline{a}' \underline{X}). \quad (*)$$

$\text{Var}(\underline{a}' \underline{X}) = \underline{a}' \Sigma \underline{a}$, so that (*) becomes:

$$\max_{\underline{a} \in \mathbb{R}^p: \|\underline{a}\|=1} \underline{a}' \Sigma \underline{a} = \max_{\underline{a} \in \mathbb{R}^p} \frac{\underline{a}' \Sigma \underline{a}}{\underline{a}' \underline{a}}, \text{ so } \left\| \frac{\underline{a}}{\|\underline{a}\|} \right\| = 1.$$

Lemma: let B $p \times p$ positive semi-definite matrix,

and $B = \sum_{i=1}^p \lambda_i \underline{e}_i \underline{e}_i^T$ its spectral decomposition. Then:

$$1) \max_{\underline{x} \in \mathbb{R}^p} \frac{\underline{x}' B \underline{x}}{\underline{x}' \underline{x}} = \lambda_1, \quad \arg \max_{\underline{x} \in \mathbb{R}^p} \frac{\underline{x}' B \underline{x}}{\underline{x}' \underline{x}} = \underline{e}_1.$$

$$2) \max_{\substack{\underline{x} \in \mathbb{R}^p: \\ \underline{x} \perp \underline{e}_1}} \frac{\underline{x}' B \underline{x}}{\underline{x}' \underline{x}} = \lambda_2, \quad \arg \max \dots = \underline{e}_2.$$

$$p) \max_{\underline{x} \in \mathbb{R}^p: \underline{x} \perp \underline{e}_2, \dots, \underline{e}_{p-1}} \frac{\underline{x}' B \underline{x}}{\underline{x}' \underline{x}} = \lambda_p = \min_{\underline{x} \in \mathbb{R}^p} \frac{\underline{x}' B \underline{x}}{\underline{x}' \underline{x}}.$$

Proof: $B = \sum_{i=1}^p \lambda_i \underline{e}_i \underline{e}_i' = P \Lambda P'$.

$$\begin{aligned} 1) \quad \frac{\underline{x}' B \underline{x}}{\underline{x}' \underline{x}} &= \frac{\underline{x}' P \Lambda P' \underline{x}}{\underline{x}' P P' \underline{x}} = \left[\text{call } \underline{y} = P' \underline{x} \right] \\ &= \frac{\underline{y}' \Lambda \underline{y}}{\underline{y}' \underline{y}} = \frac{\sum_{i=1}^p \lambda_i y_i^2}{\sum_{i=1}^p y_i^2} \leq \lambda_1 \frac{\sum_{i=1}^p y_i^2}{\sum_{i=1}^p y_i^2} = \lambda_1. \Rightarrow \lambda_1 = \max \dots \end{aligned}$$

Take now $\underline{x} = \underline{e}_1$. Then:

$$\frac{\underline{x}' B \underline{x}}{\underline{x}' \underline{x}} = \frac{\underline{e}_1' B \underline{e}_1}{\underline{e}_1' \underline{e}_1} = \lambda_1 \underline{e}_1' \underline{e}_1 = \lambda_1 \Rightarrow \underline{e}_1 = \arg \max \dots$$

$$\begin{aligned} 2) \quad \max_{\substack{\underline{x} \in \mathbb{R}^p: (\underline{x} \perp \underline{e}_1)}} \frac{\underline{x}' B \underline{x}}{\underline{x}' \underline{x}} &= \left\{ \begin{array}{l} \underline{y} = P' \underline{x} = \begin{pmatrix} \underline{e}_1' \\ \underline{e}_2' \\ \vdots \\ \underline{e}_p' \end{pmatrix} \underline{x} \stackrel{\underline{x} \perp \underline{e}_1}{=} \begin{pmatrix} 0 \\ \underline{e}_2' \underline{x} \\ \vdots \\ \underline{e}_p' \underline{x} \end{pmatrix} \underline{y} = \\ = \frac{\underline{y}' \Lambda \underline{y}}{\underline{y}' \underline{y}} = \frac{\sum_{i=2}^p \lambda_i y_i^2}{\sum_{i=2}^p y_i^2} \leq \lambda_2. \Rightarrow \max \dots = \lambda_2 \end{array} \right. \end{aligned}$$

Taking $\underline{x} = \underline{e}_2$, $\frac{\underline{x}' B \underline{x}}{\underline{x}' \underline{x}} = \frac{\underline{e}_2' B \underline{e}_2}{\underline{e}_2' \underline{e}_2} = \lambda_2 \underline{e}_2' \underline{e}_2 = \lambda_2$.
 $\Rightarrow \arg \max \dots = \underline{e}_2$.

p) Same for max, argmax.

In order to prove the minimum:

$$\frac{\underline{x}' B \underline{x}}{\underline{x}' \underline{x}} = \frac{\underline{y}' \Lambda \underline{y}}{\underline{y}' \underline{y}} = \frac{\sum_{i=1}^p \lambda_i y_i^2}{\sum y_i^2} \geq \lambda_p \rightarrow \lambda_p = \min \dots$$

$\underline{x} = \underline{e}_p \Rightarrow \underline{e}_p = \arg \min \dots$ (end)

Using the lemma, if $\sum = \sum_{i=1}^p \lambda_i \underline{e}_i \underline{e}_i'$, then:

$$\begin{cases} \max_{\underline{a} \in \mathbb{R}^p: \|\underline{a}\|=1} \text{Var}(\underline{a}' \underline{x}) = \lambda_1. \\ \arg \max_{\underline{a} \in \mathbb{R}^p: \|\underline{a}\|=1} \text{Var}(\underline{a}' \underline{x}) = \underline{e}_1. \end{cases}$$

Definition: $y_1 := \underline{e}_1' \underline{x}$ is called first principal component (PC1).

(Note: $y_1 = \underline{e}_1' (\underline{x} - \mu)$.)

Problem: $\max_{\underline{a} \in \mathbb{R}^p : \|\underline{a}\|=1} \text{Var}(\underline{a}' \underline{X}) = \max_{\underline{a} \in \mathbb{R}^p} \frac{\underline{a}' \Sigma \underline{a}}{\underline{a}' \underline{a}}.$ (*)

$$\text{Cov}(\underline{a}' \underline{X}, \underline{e}_1' \underline{X}) = 0$$

$$\text{Cov}(\underline{a}' \underline{X}, \underline{e}_i' \underline{X}) = 0$$

How can we express $\text{Cov}(\underline{a}' \underline{X}, \underline{b}' \underline{X})$?

Taking $C \in \mathbb{R}^{k \times p}$, $\text{Cov}(C \underline{X}) = C \Sigma C'$.

If $C = \begin{bmatrix} \underline{a}' \\ \underline{b}' \end{bmatrix} \in \mathbb{R}^{2 \times p}$, then: $\text{Cov}(C \underline{X}) = \text{Cov}\left(\begin{pmatrix} \underline{a}' \underline{X} \\ \underline{b}' \underline{X} \end{pmatrix}\right) = \begin{bmatrix} \underline{a}' \\ \underline{b}' \end{bmatrix} \Sigma \begin{bmatrix} \underline{a} & \underline{b} \end{bmatrix}$

$$= \begin{bmatrix} \underline{a}' \Sigma \underline{a} & \underline{a}' \Sigma \underline{b} \\ \underline{b}' \Sigma \underline{a} & \underline{b}' \Sigma \underline{b} \end{bmatrix}.$$

$$\Rightarrow \text{Cov}(\underline{a}' \underline{X}, \underline{b}' \underline{X}) = \underline{a}' \Sigma \underline{b}. \text{ Hence:}$$

$$\text{Cov}(\underline{a}' \underline{X}, \underline{e}_i' \underline{X}) = \underline{a}' \Sigma \underline{e}_i = \lambda_i \underline{a}' \underline{e}_i = 0 \iff \boxed{\underline{a} \perp \underline{e}_i}$$

Then we can reformulate (*) as:

$$\max_{\substack{\underline{a} \in \mathbb{R}^p : \\ \|\underline{a}\|=1 \\ \underline{a} \perp \underline{e}_1}} \frac{\underline{a}' \Sigma \underline{a}}{\underline{a}' \underline{a}} = \lambda_2.$$

(Lemma)

$$\text{argmax}_{\dots} \dots = \underline{e}_2.$$

Definition: $\boxed{y_2 = \underline{e}_2' \underline{X}}$ is PC2 (or second principal component).

(Note: other possibility is to take $y_2 = \underline{e}_2' (\underline{X} - \underline{\mu})$.)

We can generalize easily the solution of the problem for $j \geq 2$:

General problem:

$$\max_{\substack{\underline{a} \in \mathbb{R}^p : \\ \|\underline{a}\|=1}} \text{Var}(\underline{a}' \underline{X}) = \lambda_j, \quad \text{argmax}_{\dots} \dots = \underline{e}_j.$$

$$\|\underline{a}\|=1,$$

$$\text{Cov}(\underline{a}' \underline{X}, \underline{e}_i' \underline{X}) = 0$$

$$\forall i = 1, \dots, j-1$$

Definition: $\boxed{Y_j = \underline{e}_j' \underline{X}}$ is PC $_j$ (j -th principal component),
 $(Y_j = \underline{e}_j' (\underline{X} - \underline{\mu}))$

$\boxed{Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_p \end{pmatrix} = P' \underline{X}}$ is the vector of p -principal components.

Properties:

$$\bullet E[Y] = E[P' \underline{X}] = P' \underline{\mu}. \quad \left[\text{if } Y = P' (\underline{X} - \underline{\mu}) \Rightarrow E[Y] = \underline{0} \right]$$

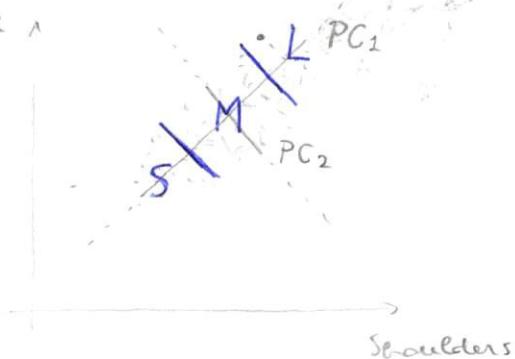
$$\bullet \text{Cov}(Y) = \text{Cov}(P' \underline{X}) = P' \Sigma P = P' P \Lambda P' P = \Lambda.$$

$$g_0 = \begin{cases} \text{cov}(Y_i, Y_j) = 0 & \text{if } i \neq j, i, j = 1 \dots p. \\ \text{Var}(Y_i) = \lambda_i. \end{cases}$$

In this way, I've constructed a powerful tool of dimension reduction.

Example: you are a tailor, and you measure size of neck & shoulders.

Neck



We can explain most of variability only with PC₁. So we project all the clients along PC₁, and divide PC₁ in three regions, {S, M, L}.

You want a better shirt? Go to a tailor! :)

Observations:

$$1) \text{General. Variance}(Y) = \text{Det}(\Lambda) = \prod_{i=1}^p \lambda_i = \text{Det}(\Sigma) = \text{G.V.}(X).$$

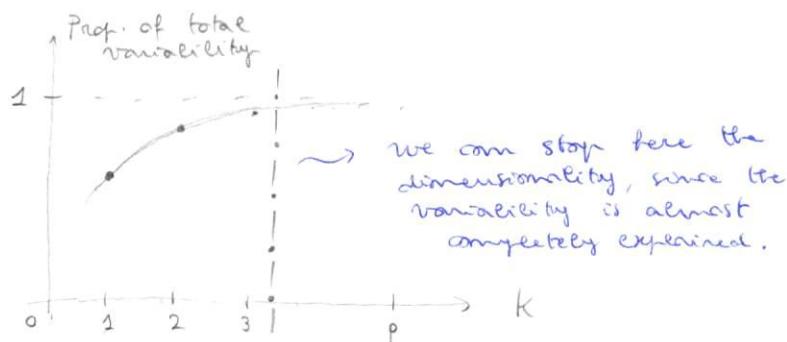
(G.V.(Y))

$$\text{Total Variance}(Y) = \text{tr}(\Lambda) = \text{tr}(\Sigma) = \text{total variance}(X).$$

We are not losing anything in the overall variance of the system.

$$2) \text{Var}(Y_1) = \lambda_1: Y_1 \text{ is capturing } \frac{\lambda_1}{\sum_i \lambda_i} \text{ proportion of total variability.}$$

$$\text{Var}(Y_2) = \lambda_2: \{Y_1, Y_2\} \text{ are capturing } \frac{\lambda_1 + \lambda_2}{\sum_i \lambda_i} \text{ of total variability.}$$



Remarks: interpretation of the PCs.

Going back to the tailor, what is the score {S, M, L} speaking about? Neck or shoulders? It's a combination of the two: it doesn't have to make sense physically, it explains the variability of data.

p-value: suppose that $\beta_1 = 0$, so H_0 hypothesis is true.

Then what is the probability of observing an estimate such $\hat{\beta}_1$?

That probability is the p-value.

If p-value is small, then what we saw is "a miracle"

if H_0 is true: H_0 is not acceptable.

Take-home message: exploring the data with graphic tools is very useful.

Amenyism: functional Principal Components.

Book: "How to Lie with Statistics".

PCA (Continued):

23/3/20

X: random vector, $E[\underline{X}] = \mu$, $Cov(\underline{X}) = \Sigma$.

$$\Sigma = \sum_{i=1}^p \lambda_i e_i e_i' \quad \text{spectral decomposition}$$

$$\lambda_1 > \lambda_2 > \dots > \lambda_p > 0.$$

$Y_i = e_i' \underline{X}$: i-th principal component.

$$Y_i = \underbrace{e_{1i} X_1 + \dots + e_{pi} X_p}_{\text{loadings}}.$$

loadings

Proposition: $\text{Cov}(Y_i, X_k) = \frac{e_{ki} \sqrt{\lambda_i}}{\sqrt{\sigma_{kk}}}, \quad i, k = 1 \dots p.$

Proof:

$$\text{Cov}(Y_i, X_k) = \frac{\text{Cov}(Y_i, X_k)}{\sqrt{\lambda_i \cdot \sigma_{kk}}}.$$

$$\text{Cov}(Y_i, X_k) = \text{Cov}(e_i' \underline{X}, u_k' \underline{X}), \quad u_k := (0 \dots 0 \underset{k}{1} 0 \dots 0).$$

$$= e_i' \sum u_k = u_k' \sum e_i =$$

$$= \lambda_i u_k' e_i = \lambda_i e_{ki}, \Rightarrow \text{Cov}(Y_i, X_k) = \frac{e_{ki} \sqrt{\lambda_i}}{\sqrt{\sigma_{kk}}} \quad (\text{end})$$

e_{ki} : k-th component of the i-th eigenvector.

Principal components obtained from standardized values:

$$\underline{z} = V^{-\frac{1}{2}} (\underline{x} - \underline{\mu}) = \left(\frac{x_1 - \mu_1}{\sqrt{\sigma_{11}}}, \dots, \frac{x_p - \mu_p}{\sqrt{\sigma_{pp}}} \right)'$$

Standardized Random Vector

$$\Rightarrow E[\underline{z}] = \underline{0}, \quad \text{cov}(\underline{z}) = V^{-\frac{1}{2}} \Sigma V^{-\frac{1}{2}} = \underline{\rho}.$$

In this way, I can get rid of the unit of measures.

We can express

$$\underline{\rho} = \sum_{i=1}^p \lambda_i \underline{e}_i \underline{e}_i'$$

where now \underline{e}_i are eigenvectors of $\underline{\rho}$ and λ_i are eigenvalues of $\underline{\rho}$, correlation matrix of \underline{X} .

The analysis proceeds as before:

$$Y_i = \underline{e}_i' \underline{z} = \underline{e}_i' V^{-\frac{1}{2}} (\underline{x} - \underline{\mu}). \quad i\text{-th principal component of } \underline{z}$$

We notice that:

$$1) \quad \sum_{i=1}^p \text{Var}(Y_i) = \text{tr}(\underline{\rho}) = p = \sum_{i=1}^p \text{Var}(z_i).$$

$$2) \quad \text{Cov}(Y_i, z_k) = \frac{\underline{e}_{ik} \sqrt{\lambda_i}}{\sqrt{1}} = \underline{e}_{ik} \sqrt{\lambda_i}.$$

3). Proportion of total variability explained by the first K components:

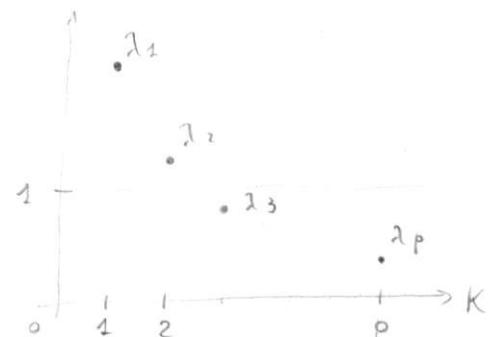
$$\left| \sum_{i=1}^K \lambda_i / p \right|. \quad K = 1, \dots, p.$$

$$4). \quad \bar{\lambda} = \frac{\sum_{i=1}^p \lambda_i}{p} = \frac{\text{tr}(\underline{\rho})}{p} = \frac{p}{p} = 1.$$

\Rightarrow The average value of the eigenvalues in $\underline{\rho}$ is 1.
This motivates us to choose Y_i if $\lambda_i > 1$.

Example: $\Sigma = \begin{pmatrix} 1 & 4 \\ 4 & 100 \end{pmatrix}$

$$\Rightarrow \underline{\rho} = \begin{pmatrix} 1 & 0.4 \\ 0.4 & 1 \end{pmatrix}.$$



$\Sigma_{22} \gg \Sigma_{11}$. So doing the PCA in Σ could not be the best thing.

For example:
 X_1 in cm.
 X_2 in mm.

In this situation, the variability of the higher range will totally mask the variability of the other variable.